



# PhaseMP: Robust 3D Pose Estimation via Phase-conditioned Human Motion Prior

Mingyi Shi<sup>1</sup>, Sebastian Starke<sup>2</sup>, Yuting Ye<sup>2</sup>, Taku Komura<sup>1</sup>, Jungdam Won<sup>3</sup>

<sup>1</sup>The University of Hong Kong, <sup>2</sup>Meta, <sup>3</sup>Seoul National University



## Motivation

- Temporal coherence plays an important role to produce realistic human motion.
- Simple interpolation is commonly used to refine motion jitters - but it doesn't work well for long-term and heavy-occluded frames.
- A periodic feature, called Phase, shows big potential to improve motion quality by describing motion in multi-dimensional sinusoidal space.

## Contribution

- We propose a **novel motion prior** based on the phase manifold for synthesizing feasible motions for various downstream tasks.
- A **new optimization framework** incorporating phase feature energy, which can work robustly for many challenging scenarios where the observation is incomplete or ambiguous in temporal and spatial domains

## Formulation

### Phase: A multi-dimensional sinusoidal vector

Given a windows of 2 seconds motion data, we take the joint velocity  $X_t \in \mathbb{R}^{3 \times J \times N}$  as input, followed by differentiable FFT layer:

$$A_t, B_t, F_t = \text{FFT}(\text{Conv}(X_t)) \quad (1)$$

Then the Phase feature  $P_t$  is defined with:

$$P_t = [p_t, F_t, A_t], \quad p_t = (A_t \cdot \sin(2\pi \cdot S_t), A_t \cdot \cos(2\pi \cdot S_t)) \quad (2)$$

### Prior: A autoregressive generation network

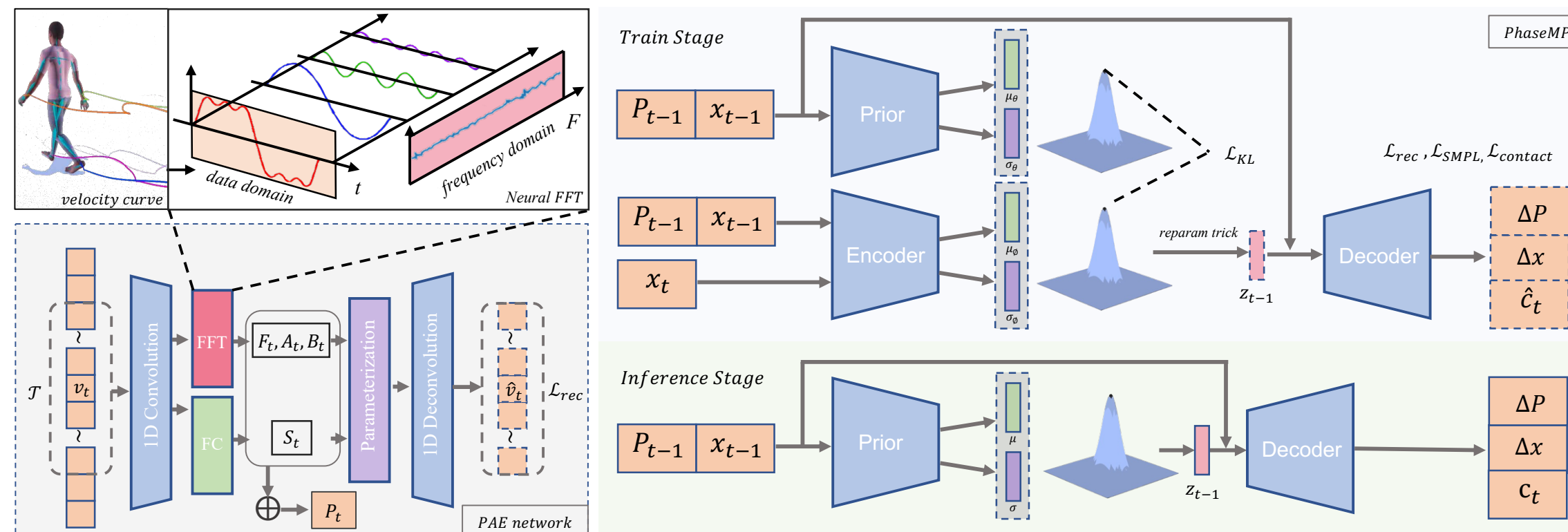
Model the transition between two frames  $[x_{t-1}, x_t]$ , with condition variant  $z$  and phase  $P$

$$\Delta x, \Delta P, c_t = G(x_{t-1}, z_{t-1}, P_{t-1}) \quad (3)$$

$$\hat{x}_t = x_{t-1} \oplus \Delta x, \hat{P}_t = P_{t-1} + \Delta P \quad (4)$$

Then a sequence of poses can be generated.

## Framework



**Left:** Phase feature extraction. The key module is a periodic auto-encoder equipped with convolution and FFT (Fast Fourier Transform) layers in its intermediate structure, allowing it to compute embeddings in the frequency domain given joint velocities as inputs.

**Right:** Conditional human motion prior. The pose in the next frame is predicted by sampling from a Gaussian distribution produced by the prior model. In the **training** stage, the prior model  $R$  is trained with posterior  $E$  by aligning their output distributions, without the input of paired frames. Thus, the prior model only predicts solely at **inference** time by only considering the previous frame. Note that we used a sine activation layer in the decoder.

## Test-time optimization: Refine the pose sequence

Given partial observations, such as a 2D landmark sequence or a partial 3D joint sequence, we estimate the original 3D pose sequence through optimization.

Stage 1: Produce initial guess of [body pose parameter  $\hat{x}_{0:T}$ , environment parameters].

Stage 2: Produce the initial [phase feature  $P_{1:T}$ , variant feature  $z_{1:T}$ ] using the encoder  $E$ . Then produce a target phase curve based on our cyclic updating and robust blending strategy.

$$p_t = \bar{A}_t \cdot I(\alpha_p)(R(\theta) \cdot p_{t-1}, (p_{t-1} + \Delta p)), \quad \theta = \Delta t \cdot 2\pi \cdot \bar{F}_t \quad (5)$$

Where the  $\bar{A}_t$  and  $\bar{F}_t$  is dynamic blended based on the confidence value from observation.

Stage 3: Refine the initial data with the following energy:

$$\text{argmin}_{z_{1:T}, \beta} (E_{obs} + E_{prior} + E_{reg} + E_{phase}) \quad (6)$$

## Experiments

We show the evaluation on following tasks: (i) motion generation (2) motion estimation from sparse observation.

Model	per-frame reconstruction			sampling (10s)		
	Contact $\uparrow$	MPJPE $\downarrow$	PJPE-std $\downarrow$	Contact $\uparrow$	ADE $\downarrow$	FDE $\downarrow$
HuMoR(MLP, w/o Phase)	0.9770	0.022	0.051	0.8216	45.43	62.47
Ours(MLP, with Phase)	0.9764	0.020	0.040	0.8525	39.48	54.96
Ours(SirenMLP, w/o Phase)	0.9788	0.019	0.031	0.8577	<b>35.47</b>	49.95
Ours(SirenMLP, with Phase)	<b>0.9799</b>	<b>0.017</b>	<b>0.021</b>	<b>0.8662</b>	42.12	<b>49.47</b>

Table 1. Comparison results on AMASS dataset reconstruction

Method	Input Conditions	fitting (3s)						
		Vis	Occ	All	Contact $\uparrow$	Accel	P-Frep	P-Dis
VPoser-t	$J_{height} > 0.9$	<b>0.67</b>	20.76	9.22	-	5.71	16.77%	2.28
MVAE		2.39	19.15	9.52	-	7.12	3.15%	0.30
HuMoR		1.46	17.40	8.24	<b>0.89</b>	5.38	3.31%	<b>0.26</b>
Ours		3.94	<b>15.63</b>	<b>8.31</b>	<b>0.89</b>	<b>4.58</b>	<b>3.04%</b>	0.28
HuMoR	$J_{end\ effectors}$	<b>3.05</b>	4.12	3.83	0.96	4.91	0.31%	1.03
Ours		3.16	<b>4.07</b>	<b>3.79</b>	<b>0.97</b>	<b>4.88</b>	<b>0.28%</b>	<b>1.02</b>
HuMoR	10 frames interval	5.56	7.76	7.49	0.91	7.72	1.57%	1.90
Ours		<b>3.19</b>	<b>4.92</b>	<b>4.33</b>	<b>0.93</b>	<b>6.33</b>	<b>1.31%</b>	<b>1.72</b>

Table 2. Comparison results on estimation from different input conditions

## Results

